# EXHIBIT 155

Page 1

UNITED STATES DISTRICT COURT

NORTHERN DISTRICT OF CALIFORNIA

SAN FRANCISCO DIVISION

-------------------------------------

RICHARD KADREY, ET AL.,                )

    Individual and Representative )

                       Plaintiffs,  ) Lead Case No.

             v.                   ) 3:23-cv-03417-VC

META PLATFORMS, INC.,                  )

                Defendant.   )

-------------------------------------

* * * HIGHLY CONFIDENTIAL * * *

* * * ATTORNEYS' EYES ONLY * * *

VIDEO-RECORDED 30(b)(6) DEPOSITION OF

MICHAEL CLARK (mitigation)

MONDAY, MARCH 3, 2025

DENVER, COLORADO

2:59 P.M. MST

REPORTED BY KATHY L. DAVIS, CRR, CMR

-----------------------------------------------------

DIGITAL EVIDENCE GROUP

1730 M. street, NW, Suite 812

Washington, D.C. 20036

(202) 232-0646

Page 43

```
 1   repetitive and duplicative."  Repetitive was the

 2   first bullet.

 3        Q     Mm-hmm.

 4        A     And then duplicative and repetitive are

 5   kind of the same.

 6        Q     Mm-hmm.

 7        A     My phrase in my prior testimony feels

 8   incomplete, but please ask your question.

 9        Q     So for the examples of filter data and

10   repetition, he provides copyright information as

11   repetitive data, he provides personal identifying

12   information as repetitive data, right?

13        A     We did miss a few bullets, Remove

14   excessive new line crack terms and Remove documents

15   in the token distribution tool.

16        Q     Mm-hmm.

17        A     Um -- I'm still unclear on the question.

18   I apologize.

19        Q     My question is that when you're talking

20   about this script that he wrote to remove data that

21   had a higher likelihood of being regurgitated, among

22   the repetitive data that he removes is also --

23   it's -- is copyright data; is that correct?

24        A     Um, among the data that's being removed,

25   the repetitive and duplicative data has a higher
```

Page 44

```
 1   likelihood of being regurgitated.  The copyright

 2   information and line spacing and emails, part of that

 3   is for preventing regurgitation, but also just

 4   general stability of the model from -- from a

 5   language perspective.

 6        Q     Okay.

 7              MS. POUEYMIROU:  Can we turn to Tab 15.

 8   This will be 14, Exhibit 14.

 9              (Clark Exhibit 14, marked for

10              identification.)

11        Q     (BY MS. POUEYMIROU)  Are you familiar

12   with this document, Mr. Clark?

13        A     I don't remember this document.

14        Q     Are you familiar with Meta's source

15   metadata project?

16        A     Let me read this to see if I am.

17              MS. POUEYMIROU:  How much time?

18              VIDEOGRAPHER:  50 minutes.

19              MS. POUEYMIROU:  5-0?

20              VIDEOGRAPHER:  (Nodded.)

21              MS. POUEYMIROU:  Thought it was 10.

22   Almost gave me a heart attack.

23              (A discussion was had off the record.)

24        A     Okay.

25        Q     (BY MS. POUEYMIROU)  Okay.  So the
```

Page 45

```
 1  concern about source metadata leaking leads to Meta

 2  employees talking about whether it can be hashed or

 3  unhashed in use cases; is that correct?

 4       A    More than just source metadata leaking as

 5  potential harm.  What's -- what's the question?

 6       Q    My question is, they're trying to figure

 7  out a safe way to train and source metadata without

 8  risking IP or brand safety.  Is that -- is that fair

 9  to say?

10            MR. WEINSTEIN:  Object to form.

11       A    With the -- going back to Exhibit 13, it

12  is much more than just leakage of info.  It is also

13  training stability and performance, that what does

14  either the hash or the raw metadata introduce from

15  a -- from a instability perspective into the model?

16       Q    (BY MS. POUEYMIROU)  Okay.  I recognize

17  that, but another consideration is about IP and brand

18  safety; is that correct?

19       A    That is one of the three that is listed.

20       Q    Okay.  And so the document that you now

21  have before you -- this is Exhibit 14 -- is entitled

22  Source Metadata Strategy.  Do you see that?

23       A    I do see that.

24       Q    And I had some questions about that.

25  First, are you familiar with this source metadata
```

Page 46

```
 1   project at Meta?

 2        A       Personally, I was familiar that this was

 3   a -- this was work that was being evaluated.

 4        Q       And why were you personally familiar with

 5   that?

 6        A       Just in -- working in generative AI and

 7   seeing other -- other projects that were in flight is

 8   how I'm familiar -- or were in the works is why I was

 9   familiar with it.

10        Q       And how did you understand the project?

11        A       I understood the project as, Could adding

12   source metadata to the model help validate where we

13   see issues or regression in certain types of

14   performance, know where that came from, but also from

15   an academic perspective in trying to get to a place

16   to where as -- the question of, Was Abraham Lincoln a

17   president, where that knowledge came from.

18        Q       So the -- the metadata that Nikolay

19   Bashlykov had stripped from LibGen, for example, was

20   metadata that Meta is now considering how it can

21   safely train on; is that correct?

22            MR. WEINSTEIN:  Object to form.

23        A       The data that Mr. Bashlykov's script

24   removed was a subset of the overall data that was

25   being tested and evaluated as part of the source
```

Page 47

 1   metadata strategy.  However, in looking and comparing

 2   this data, the data in Mr. Bashlykov's were where

 3   lines contained that --

 4              (Weather noise was heard.)

 5              THE DEPONENT:  Is that all wind?

 6              MS. POUEYMIROU:  It sounds like it.

 7              MR. WEINSTEIN:  Okay.

 8       A    Okay.  Apologies.  To get back to

 9   answering the question, what Mr. Bashlykov stripped

10   out were individual lines --

11       Q    (BY MS. POUEYMIROU)  Mm-hmm.

12       A    -- but not actually the full set of

13   metadata or source metadata that would have been

14   needed to have filled in -- and if I look at the

15   example, Title, Type, Host Name, URL, Date,

16   Generator, Source, Media --

17       Q    Mm-hmm.

18       A    -- Media Bias, and other information to

19   fill out the metadata fields.

20       Q    So what's the URL when you are -- you've

21   gotten your data from LibGen?  Will that be in the

22   URL?

23       A    This was -- the section I just read was

24   specific to Common Crawl.

25       Q    Okay.  And so if you go to the next page,

Page 53

 1    documents from other URLs?"  What is this concern

 2    about people figuring out the hash associated with

 3    Wikipedia?

 4                MR. WEINSTEIN:  Object to form.

 5        A    I would have to speculate that the harm

 6    is going back to releasing information about the data

 7    mix that was used to build out the model.  I believe

 8    that comment was from Eric Smith, not Soumya.

 9        Q    (BY MS. POUEYMIROU)  Oh, it looked like

10    it was on the back.

11        A    Oh, I was on 81207.00005.

12        Q    Yeah.

13        A    Okay.  Eric Smith had made that comment,

14    just for accuracy.

15        Q    Oh, okay.

16                So the leakage of source metadata

17    exposing to the public what Meta had trained on --

18                MR. WEINSTEIN:  Object to form.

19        A    The proprietary data mix that was used

20    for developing the model.

21        Q    (BY MS. POUEYMIROU)  And what do you mean

22    by "proprietary data mix"?

23        A    The combination of data that was

24    developed -- the combination of data -- the exact

25    combination of data that was used to train the model.

Page 54

```
 1        Q      How would -- how would a person know the

 2   combination or the data mix of your data --

 3               MR. WEINSTEIN:  Object to form.

 4        A      I --

 5        Q      (BY MS. POUEYMIROU)  -- as opposed to

 6   just data that was trained on?

 7               MR. WEINSTEIN:  Object to form.

 8        Q      (BY MS. POUEYMIROU)  I have always

 9   thought of data mix in this case as how much of one

10   data mix compared -- one dataset compared to another.

11   Are you just talking about the data that Meta trained

12   on being released to the public?

13        A      The com- -- the exact combination of what

14   datasets and/or data Meta trained on for that model.

15        Q      How does the -- how does the training on

16   the title and author of a work --

17               MR. WEINSTEIN:  Object to form.

18               MS. POUEYMIROU:  I didn't finish.

19               MR. WEINSTEIN:  Oh, sorry.  Sorry about

20   that.

21        Q      (BY MS. POUEYMIROU)  -- expose a data

22   mix?

23               MR. WEINSTEIN:  Object to form, scope.

24        A      The -- only the -- as the motivation

25   said, the reason that you would train on it is to
```

Page 55

```
 1   make it available from an explainability or to

 2   improve the quality from a factuality perspective.

 3              And so there's no sense in training -- in

 4   adding additional data that can't be used, which

 5   therefore exposes that data.

 6       Q     (BY MS. POUEYMIROU)  So I'm trying to

 7   understand, though, your point about proprietary

 8   datasets.  Are you saying that LibGen is Meta's

 9   intellectual property and if training on source

10   metadata from LibGen, which then gets regurgitated,

11   it would expose?

12              MR. WEINSTEIN:  Object to form.

13       A     Just to -- I believe that's a

14   misclarification of the phrase that I used.  It's the

15   proprietary data mix, the combination of what data

16   was used to train not the data itself being

17   proprietary.

18       Q     (BY MS. POUEYMIROU)  So you're saying the

19   fact that Meta trained on Common Crawl and Wikipedia

20   and Substack and Archive and . . .  You're saying

21   that is the concern, that that would be released to

22   the public?

23              MR. WEINSTEIN:  Object to form.

24       A     That would be a -- that is a potential

25   example, yes.
```

Page 56

```
 1        Q      (BY MS. POUEYMIROU)  But didn't Meta

 2   speak openly about its use of Common Crawl and other

 3   datasets in its published papers?

 4              MR. WEINSTEIN:  Object to form.  Outside

 5   the scope.

 6        A      It only spoke about that in LLaMA 1,

 7   which was a research release only, and not in any

 8   other models after.  That was the one mitigation that

 9   I forgot in my list when I read them off, was that --

10   to not share that data mix, to not share the data mix

11   as part of one of the mitigations.

12        Q      (BY MS. POUEYMIROU)  Okay.  And so your

13   testimony is that by removing the source metadata,

14   Meta was able to conceal what it trained on because

15   by training on source metadata, the regurgitation

16   risk was higher?

17              MR. WEINSTEIN:  Object to form.

18        A      I believe that's a mischaracterization of

19   my testimony as the source metadata document that was

20   listed here was adding additional data into the

21   training mix and not taking it out.

22        Q      (BY MS. POUEYMIROU)  Was that --

23        A      And --

24        Q      -- because it had first been taken out?

25              MR. WEINSTEIN:  Object to form.
```

Page 57

```
 1       A      This -- sorry.  I wasn't complete with my

 2   answer.  Do you mind asking the original question?

 3       Q      (BY MS. POUEYMIROU)  Sure.  The decision

 4   to add source metadata to training is predicated on

 5   the fact that Meta had already removed source

 6   metadata --

 7              MR. WEINSTEIN:  Object to form.

 8       Q      (BY MS. POUEYMIROU)  -- from training; is

 9   that right?

10              MR. WEINSTEIN:  Sorry.  Object to form.

11   Lacks foundation.

12       A      If I look at the examples from -- which

13   were in 65250 from what Mr. Bashlykov had removed,

14   this information is incomplete in how it's structured

15   and was not trained in a structured way that could

16   have been used in the same way as the source metadata

17   strategy.  This was removing specific things where

18   they were placed in the doc that were not structured

19   or in a way that could actually help what the

20   motivation of the test that was being done around the

21   source metadata strategy was, which, to the best of

22   my knowledge, was not actually pursued because of

23   challenges it created in both the performance and

24   stability of the model when it was added.

25       Q      (BY MS. POUEYMIROU)  And also with
```

Page 58

 1    respect to --

 2         A     The reasons it was not chosen is it

 3    created too much regression, added too much

 4    additional volume, which you can see reflected in the

 5    comments of data that wasn't actually usable as part

 6    of inference and was therefore not pursued because of

 7    the challenges around performance and stability of

 8    the model.

 9         Q     You're saying that with respect to

10    personal identifying information too?  That was why

11    it was removed?

12              MR. WEINSTEIN:  Object to form.

13         A     We were just talking about the source

14    metadata strategy project, which was explicitly

15    adding multiple structured fields into the model.

16    The PII that was removed as part of Mr. Bashlykov's

17    script was to prevent repetitive personal information

18    from being able to be regurgitated about private

19    individuals.  That might have been email or other

20    kinds of PII.

21         Q     (BY MS. POUEYMIROU)  And the use of that

22    script to remove the copyright information was also a

23    regurgitation concern?

24              MR. WEINSTEIN:  Object to form.

25         A     That was both regurgitation, plus the

Page 59

```
 1   repetitive characters of copyright and the C -- the

 2   copyright mark -- copymark would affect the

 3   performance of the model as well because that would

 4   be very much repetitive and become part of outputted

 5   responses, but not necessarily adjacent to

 6   copyrighted information.

 7          Q       (BY MS. POUEYMIROU)  Okay.

 8                  MS. POUEYMIROU:  Can we pull out 23.

 9                  (Clark Exhibit 15, marked for

10                  identification.)

11          A       Are we switching topics?

12          Q       (BY MS. POUEYMIROU)  We're switching

13   documents.

14          A       Okay.

15          Q       The topic is mitigation, though.  Now,

16   this is a multipage document.  I'm only going to be

17   asking about parts later on, so just take your time.

18                  MS. POUEYMIROU:  This is Exhibit --

19                  THE REPORTER:  15.

20          Q       (BY MS. POUEYMIROU)  -- 15.  All right.

21   I want to look at the page ending in 2225.  Do you

22   want to take a look at that page?

23                  MS. POUEYMIROU:  How much time do we

24   have?

25                  VIDEOGRAPHER:  24 minutes.
```

Page 75

1    STATE OF COLORADO     )

2                          )ss.    REPORTER'S CERTIFICATE

3    COUNTY OF DENVER      )

4         I, Kathy L. Davis, do hereby certify that I am a

5    Registered Professional Reporter within the State of

6    Colorado; that previous to the commencement of the

7    examination, the deponent was duly sworn to testify

8    to the truth.

9         I further certify that this deposition was taken

10   in shorthand by me at the time and place herein set

11   forth, that it was thereafter reduced to typewritten

12   form, and that the foregoing constitutes a true and

13   correct transcript.

14        I further certify that I am not related to,

15   employed by, nor of counsel for any of the parties or

16   attorneys herein, nor otherwise interested in the

17   result of the within action.

18        In witness whereof, I have affixed my signature

19   this 4th day of March, 2025.

20

21

22                    _Kathy L. Davis_

23                 Kathy L. Davis
                Certified Realtime Reporter

24

25